

ELM—the database of eukaryotic linear motifs

Holger Dinkel¹, Sushama Michael¹, Robert J. Weatheritt¹, Norman E. Davey¹, Kim Van Roey¹, Brigitte Altenberg¹, Grischa Toedt¹, Bora Uyar¹, Markus Seiler¹, Aidan Budd¹, Lisa Jödicke¹, Marcel A. Dammert¹, Christian Schroeter¹, Maria Hammer¹, Tobias Schmidt¹, Peter Jehl¹, Caroline McGuigan¹, Magdalena Dymecka², Claudia Chica³, Katja Luck⁴, Allegra Via⁵, Andrew Chatr-aryamontri⁶, Niall Haslam⁷, Gleb Grebnev⁷, Richard J. Edwards⁸, Michel O. Steinmetz⁹, Heike Meiselbach¹⁰, Francesca Diella^{1,11} and Toby J. Gibson^{1,*}

¹Structural and Computational Biology, European Molecular Biology Laboratory, Heidelberg, Germany,

²Laboratory of Bioinformatics and Systems Biology, M. Skłodowska-Curie Cancer Center and Institute of Oncology, WK Roentgena 5, 02-781 Warsaw, Poland, ³Genoscope (CEA – Institut de Génomique), 2 rue Gaston Cremieux CP5706, 91057 Evry, ⁴Group Oncoproteins, Unité CNRS-UDS UMR 7242, Institut de Recherche de l'Ecole de Biotechnologie de Strasbourg, 1, Bd Sébastien Brant, BP 10413, 67412 Illkirch – Cedex, France, ⁵Biocomputing Group, Department of Physics, Sapienza University of Rome, P.le Aldo Moro 5, Rome, Italy, ⁶School of Biological Sciences, University of Edinburgh, Mayfield Road, Edinburgh EH9 3JR, UK, ⁷School of Medicine and Medical Science, University College, Dublin, Ireland, ⁸Centre for Biological Sciences, Institute for Life Sciences, University of Southampton, UK, ⁹Biomolecular Research, Paul Scherrer Institut, CH-5232 Villigen PSI, Switzerland, ¹⁰Bioinformatik, Institut für Biochemie, Friedrich-Alexander-Universität, Fahrstraße 17, 91054 Erlangen-Nürnberg and ¹¹Molecular Health GmbH Belfortstr. 2, 69115 Heidelberg, Germany

Received September 13, 2011; Revised and Accepted October 27, 2011

ABSTRACT

Linear motifs are short, evolutionarily plastic components of regulatory proteins and provide low-affinity interaction interfaces. These compact modules play central roles in mediating every aspect of the regulatory functionality of the cell. They are particularly prominent in mediating cell signaling, controlling protein turnover and directing protein localization. Given their importance, our understanding of motifs is surprisingly limited, largely as a result of the difficulty of discovery, both experimentally and computationally. The Eukaryotic Linear Motif (ELM) resource at <http://elm.eu.org> provides the biological community with a comprehensive database of known experimentally validated motifs, and an exploratory tool to discover putative linear motifs in user-submitted protein sequences. The current update of the ELM database comprises 1800 annotated motif instances representing 170 distinct functional classes, including approximately 500 novel instances and

24 novel classes. Several older motif class entries have been also revisited, improving annotation and adding novel instances. Furthermore, addition of full-text search capabilities, an enhanced interface and simplified batch download has improved the overall accessibility of the ELM data. The motif discovery portion of the ELM resource has added conservation, and structural attributes have been incorporated to aid users to discriminate biologically relevant motifs from stochastically occurring non-functional instances.

INTRODUCTION

Short linear motifs (SLiMs, LMs or MiniMotifs) are regulatory protein modules characterized by their compact interaction interfaces (the affinity and specificity determining residues are usually encoded between 3 and 11 contiguous amino acids (1)) and their enrichment in natively unstructured, or disordered, regions of proteins (2). As a result of limited intermolecular contacts with their interaction partners, SLiMs bind with relatively

*To whom correspondence should be addressed. Tel: +49 (0) 6221 3878398; Fax: +49 (0) 6221 387517; Email: gibson@embl-heidelberg.de

low affinity (in the low-micromolar range), an advantageous attribute for use as transient, conditional and tunable interactions necessary for many regulatory processes. Due to the limited number of mutations necessary for the genesis of a novel motif, SLiMs are amenable to convergent evolution, functioning as a driver of network evolution by adding novel interaction interfaces, and thereby new functionality, to proteins. This evolutionary plasticity facilitates the rapid proliferation within a proteome, and as a result, motif use is ubiquitous in higher eukaryotes.

SLiMs play an important role for many regulatory processes such as signal transduction, protein trafficking and post-translational modification (3,4). Their importance to the correct functionality of the cell is also reflected by the outcome of motif deregulation. For example, point mutations in SLiMs have been shown to lead severe pathologies such as 'Noonan-like syndrome' (5), 'Liddle's syndrome' (6) or 'Retinitis pigmentosa' (7). Furthermore, mimicry of linear motifs by viruses to hijack their hosts' existing cellular machinery plays an important role in many viral life cycles (8). However, despite their obvious importance to eukaryotic cell regulation, our understanding of SLiM biology is relatively limited, and it has been suggested that, to date, we have only discovered a small portion of the human motifs (9).

Several resources are devoted to the annotation and/or detection of SLiMs [Prosite (10), MiniMotifMiner (11) and Scansite (12)]. Here, we report on the 2012 status of the Eukaryotic Linear Motif database.

THE ELM RESOURCE

The ELM initiative (<http://elm.eu.org>) has focused on gathering, storing and providing information about short linear motifs since 2003. It was established as the first manually annotated collection of SLiM classes and as a tool for discovering linear motif instances in proteins (13). As it was mainly focused on the eukaryotic sequences, it was termed the Eukaryotic Linear Motif resource, usually shortened to ELM. The ELM resource consists of two applications: the ELM database of curated motif classes and instances, and the motif detection pipeline to detect putative SLiM instances in query sequences. In the ELM database, SLiMs are annotated as 'ELM classes', divided into four 'types': cleavage

sites (CLV), ligand binding sites (LIG), sites of post-translational modification (MOD) and subcellular targeting sites (TRG) (Table 1). Currently, the ELM database contains 170 linear motif classes with more than 1800 motif instances linked to more than 1500 literature references (Table 1). Each class is described by a regular expression capturing the key specificity and affinity determining amino acid residues. A regular expression is a computer-readable term for sequence annotation and is used by the ELM motif detection pipeline to scan proteins for putative instances of annotated ELM classes. The search form for sequence input is shown in Figure 1, while the results page showing the putative and annotated instances is illustrated in Figure 2.

The ELM resource is powered by a PostgreSQL relational database for data storage and a PYTHON web framework for data retrieval/visualization. The main tables within the database contain information about ELM classes, ELM instances, sequences, references, taxonomy and links to other databases [the database structure is described in greater detail in (14)].

New ELM classes

Since the last release (14), 24 new ELM classes have been added to the ELM database (Table 1) and several more have been updated. One of the newly annotated motif classes is the AGC kinase docking motif (LIG_AGCK_PIF), consisting of three distinct classes. It is present in the non-catalytic C-terminal tail of AGC kinases that constitute a family of serine/threonine kinases consisting of 60 members that regulate critical processes, including cell growth and survival. Deregulation of these enzymes is a causative factor in different diseases such as cancer and diabetes. The motif interacts with the PDK1 Interacting Fragment (PIF) pocket in the kinase domain of AGC kinases. It mediates intramolecular binding to the PIF pocket, serving as a *cis*-activating module together with other regulatory sequences in the C-tail. Interestingly, in some kinases the motif also acts as a PDK1 docking site that *trans*-activates PDK1, which itself lacks the regulatory C-tail, by interacting with the PDK1 PIF pocket. PDK1 in turn will phosphorylate and activate the docked kinase. Other novel classes (Table 2) include phosphodegrons, which are important mediators of phosphorylation-dependent protein destruction, and the LYPxL motif, which is involved in endosomal

Table 1. Summary of data stored in the ELM database^a

Number of functional site entries	ELM motif classes	ELM motif instances	Links to PDB structures	GO terms	Pubmed links	
Totals	115	170	1840	195	340	
By category	LIG MOD TRG CLV	111 30 21 8	Human Mouse Rat Fly Yeast Other	1004 160 102 67 90 417	Biological process Cell compartment Molecular function	787 173 74 93
					From ELM motif From instance	1071

^aAs of October 2011.

ELM The Eukaryotic Linear Motif resource for *Functional Sites in Proteins*

Search ELMs Instances Candidates Links About News Help Diseases Viruses

Functional site prediction

Protein sequence

Enter Uniprot identifier or accession number: (auto-completion)
EPN1_HUMAN e.g. EPN1_HUMAN, P04637, TAU_HUMAN

Or paste the sequence (Single letter code sequence only or [FASTA format](#)):

```
>sp|Q9Y6I3|EPN1_HUMAN Epsin-1 OS=Homo sapiens GN=EPN1 PE=1 SV=2
MSTSSLRQMKNIIVNYSEAEIKVREATSNDPGPSSSLMSEIADLTYNVVAFSEIMSMI
WKRLLNDIGKNNWRHVVKAMTLMLEYLTKGSERVSQQCKENMYAVOTLKDFOYVDRDGKDOG
VNVREREAKQVQLVALLRDEDRLLREERAHALKTKEKLQAQTATASSAAVGSGPPPEAEQAWPQS
SGEEELQLQALAMSKEEADOPPSCGFEDDAQLQALSLSREEHIDKEERIRRGGDDRLRQM
ATEESKRETGGKEESSLMDLADVFTAPAPAPTIDPWGGPAPMAAAVPTAAPTSDPWGGP
VPPAADPWNGGPAPPTPASGDPMRPAAPAGPSVPDWGGTPAPAAEGEPTPDPWGSSDGVPV
SGPSASDPWTTPAPAEFDWGGSPAKPSTNGTTAAAGGFDTEDPESDFDRLRTALPTSGSS
AGELELLAGEVPARSPGAFDMSGVRGSLAEAVGSPPAAATPTPTPPTRKTPESFLGPNAA
LVQDLSLVSRPGPTPPGAKASNPFLPGGGPATGPSVTMPFOPAPPATLTNQLRLSPVPEP
VPGAPPTYISPLGGGPGLPMMPGGPAPNTNPFLI
```

Cell compartment (one or several): not specified, extracellular, nucleus, cytosol, peroxisome, glycosome, glyoxisome, Golgi apparatus, endoplasmic reticulum, lysosome, endosome, plasma membrane, mitochondrion

Context information: Type in species name (auto-completion): Homo sapiens

Motif Probability Cutoff (beta): 0.01

Submit Reset Form

Disclaimer

Short patterns applied to proteins are usually not statistically significant: Therefore we can't provide E-values as with BLAST searches. This means that most matches shown are more likely to be false positives than true matches. We hope that ELM server results will prove useful as guides to experimentation but they should not be treated as factual findings.

Feedback

If you run into any bugs, please let us know: bugs@elm.eu.org.

Comments or suggestions: feedback@elm.eu.org

Figure 1. ELM start page. The user can submit a query sequence to the motif detection pipeline either as UniProt accession number or in FASTA format. Filtering criteria such as taxonomic range or cellular compartment should be activated to limit the resulting list of SLiM instances.

sorting of membrane proteins but is also implicated in retrovirus budding.

New ELM instances

Annotated ELM instances serve as representative examples of the respective ELM class. They are also

invaluable for the computational analysis and classification of motifs (15). Therefore, special emphasis has been put on the curation of more than 500 novel ELM instances (in 40 different classes) by scanning and annotating more than 400 articles. The number of protein databank (PDB) entries annotated have been increased to 195 (Table 1), meaning that for ~10% of all instances there is a 3D

■ Summary of features reported by the ELM resource.

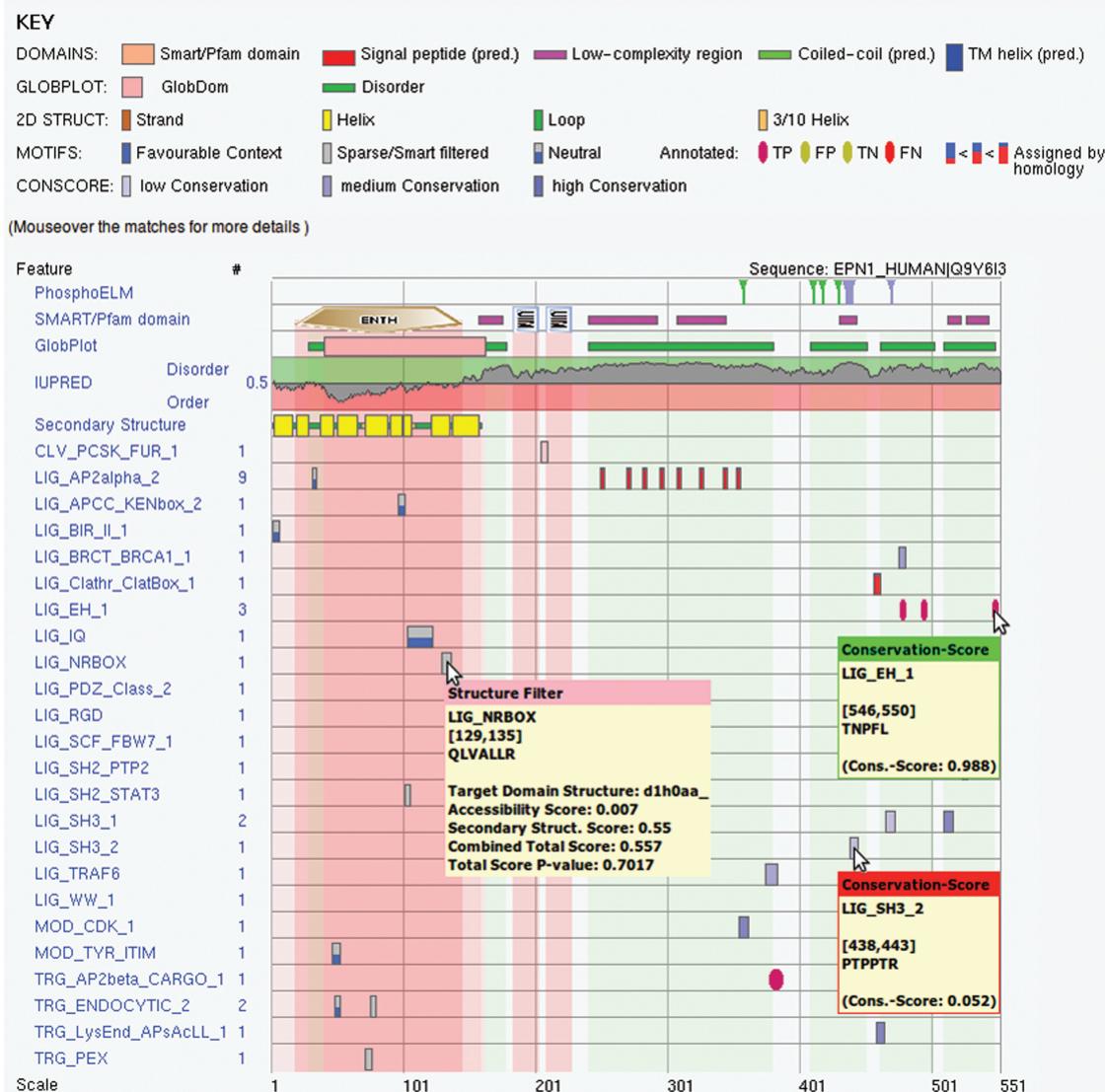


Figure 2. ELM motif detection pipeline output page. The top legend explains the different colors/symbols used. The graphical output of ELM concentrates the output of multiple sequence classification algorithms; phosphorylation sites from Phospho.ELM, protein domains detected by SMART/Pfam, disorder predictions by GlobPlot and IUPred and secondary structure (18). The lower part contains the annotated and putative ELM instances for the given protein sequence (Epsin1, UniProt accession Q9Y6I3). The background is colored according to the structural information available. Each box represents one ELM instance, the color of which indicates the likelihood that this instance is functional: grey instances are buried within structured regions, while shades of blue represent instances outside of structured regions and hint on sequence conservation, with pale blue representing weak sequence conservation and dark blue indicating strong sequence conservation. Red ellipses or boxes mark instances that are annotated in the query sequence or a homologous sequence, respectively.

protein structure annotated, giving more detailed information about the biological context of the respective motif.

NEW FEATURES

The ELM website at <http://elm.eu.org> can be used in two ways: first, as a front-end to explore the ELM database of curated ELM classes and instances, and second, to run the motif detection pipeline to detect putative SLiM instances in query sequences. Both interfaces have been improved with the most notable changes listed below.

User interface

The database user interface, having been stable for many years, has been overhauled and replaced by a novel interface introducing several new features (Figure 1). Up-to-date web technologies have been used to improve the general user experience: the PYTHON framework DJANGO (<http://www.djangoproject.com>) dynamically creates and serves all HTML pages, while JavaScript was used to make the whole site more interactive and thus improve the user experience. In particular, the ELM detail pages (Figure 3), which hold the most

Table 2. List of novel ELM classes^a

Identifier	Description
LIG_Actin_WH2_1	Motifs, present in proteins in several repeats, which mediate binding to the hydrophobic cleft created by subdomains 1 and 3 of G-actin
LIG_Actin_WH2_2	
LIG_Actin_RPEL_3	
LIG_AGCK_PIF_1	The AGCK docking motif mediates intramolecular interactions to the PDK1 Interacting Fragment (PIF) pocket, serving as a <i>cis</i> -activating module
LIG_AGCK_PIF_2	
LIG_AGCK_PIF_3	
LIG_BIR_II_1	IAP-binding motifs are found in pro-apoptotic proteins and function in the abrogation of caspase inhibition by inhibitor of apoptosis proteins in apoptotic cells
LIG_BIR_III_1	
LIG_BIR_III_2	
LIG_BIR_III_3	
LIG_BIR_III_4	
LIG_eIF4E_1	Motif binding to the dorsal surface of eIF4E
LIG_eIF4E_2	
LIG_EVH1_3	A proline-rich motif binding to EVH1/WH1 domains of WASP and N-WASP proteins
LIG_HCF-1_HBM_1	The DHxY Host Cell Factor-1 binding motif interacts with the N-terminal kelch propeller domain of the cell cycle regulator HCF-1
LIG_Integrin_isoDGR_1	Present in proteins of extracellular matrix which upon deamidation forms biologically active isoDGR motif which binds to various members of integrin family
LIG_LYPXL_L_2	The LYPxL motif binds the V-domain of Alix, a protein involved in endosomal sorting
LIG_LYPXL_S_1	
LIG_PAM2_1	Peptide ligand motif that directly interacts with the MLLE/PABC domain found in poly(A) binding proteins and HYD E3 ubiquitin ligases
LIG_PIKK_1	Motif located in the C terminus of Nbs1 and its homologous interacting with PIKK family members
LIG_Rb_pABgroove_1	The LxxLFD motif binds in a deep groove between pocket A and pocket B of the Retinoblastoma protein
LIG_SCF_FBW7_1	The TPxxS phospho-dependent degron binds the FBW7 F box proteins of the SCF (Skp1-Cullin-Fbox) complex
LIG_SCF_FBW7_2	
LIG_SPAK-OSR1_1	SPAK/OSR1 kinase binding motif acts as a docking site which aids the interaction with their binding partners including the upstream activators and the phosphorylated substrates

^aAs of October 2011.

important information about each ELM class including references, regular expression, taxonomic distribution and gene ontology terms (**Table 3**), have been updated by annotating the protein domain interacting with the respective motif. Where available, a 3D model of representative protein databank structures of linear motif interactions was added to the ELM detail page (**Figure 3**, top right).

To cope with the increasing amount of annotated classes as well as instances, a novel query interface was introduced to assist the user in finding information of interest. The ELM browser (**Figure 4**) now features a search interface for free text search. In addition, the search results can also be filtered and reordered using buttons (**Figure 4**, left side) and table headers, respectively, and be downloaded as tab-separated values (TSV).

Further, improvements to the ELM database include revising the experimental methods used for annotation by using a standardized methods vocabulary [in sync with PSI-MI ontology (16,17)].

A candidate page has been introduced to display novel ELM classes that have not yet been annotated in detail or are currently undergoing annotation. We invite researchers to send us their feedback and expert opinion on these classes and to contribute novel motif classes that will be added to the candidate page and ultimately be turned into full ELM classes (**Figure 5**). Minimum requirements are at

least one literature reference as well as a short description. In addition, a draft regular expression or a 3D structure showing the relevant interaction would also be helpful. Currently, the number of possible ELM classes on this candidate list (awaiting further annotation) exceeds the number of completely annotated classes, indicating the great demand for further annotation.

Graphical representation of sequence search

The ELM motif detection pipeline scans protein sequences for matches to the regular expressions of annotated ELM classes (**Figure 2**). The query output combines these putative instances with information from the database (annotated ELM instances) as well as predictions from different algorithms/filters. The ELM resource employs a structural filter (18) to highlight and mask secondary structure elements, as well as SMART (19) to detect protein domains. Furthermore, an additional disorder prediction algorithm (IUPred) (20) has been included to predict ordered/disordered regions within the protein. IUPred uses a cutoff of 0.5 to classify a sequence region as either structured or disordered, with values above this threshold corresponding to disorder, highlighted in green background and lower values indicating structured regions, displayed in red background in the output graph. Disorder and domain information is combined by

ELM The Eukaryotic Linear Motif resource for Functional Sites in Proteins

Search ELMs Instances Candidates Links About News Help Diseases Viruses

TRG_AP2beta_CARGO_1

<< MOD_WntLipid << Menu >> TRG_Cilium_Arf4_1 >>

Functional Site Class: AP-2 beta2 appendage CCV component motifs

Functional site description: Several motifs are responsible for the binding of accessory endocytic proteins to the beta2-subunit appendage of the adaptor protein complex AP-2 as part of their recruitment to the site of clathrin coated vesicle (CCV) formation. Proteins binding the platform subdomain have been found to be cargo family specific (for example can load all GPCRs, or all LDL receptor family members) clathrin adaptors. Accessory proteins which help in CCV formation bind the sandwich subdomain site or the alpha ear domain.

ELMs: TRG_AP2beta_CARGO_1

Description: Motif binding as a helix in a depression on the top surface of the AP-2 beta appendage platform subdomain. The pattern [ED]x(1,2)Fxx[FL]xxxR is conserved in beta Arrestins, ARH and Epsin-1, -2 of vertebrates. It is also found in homologues of other metazoans, but the pattern is sometimes not matched exactly, meaning that the ELM regular expression will not provide a match. In other lineages, if there is an equivalent motif, the pattern is likely to have diverged.

Pattern: [DE] . {1,2}F[^P] [^P][FL][^P][^P]

Present in taxons: Metazoa (Probability: 0.0000182)

PDB Structure: 2IV8

Interaction Domain:

B2-adapt-app_C (PF09066)
Beta2-adaptin appendage, C-terminal sub-domain
(Stoichiometry: 1 : 1)

See 4 Instances for TRG_AP2beta_CARGO_1

Abstract

At least two different surfaces of the AP-2 beta2 appendage domain can bind linear motifs in other endocytic regulatory proteins. The platform subdomain or top surface binds a helical [ED]x(1,2)Fxx[FL]xxxR motif found in Epsin-1 and -2 which bind ubiquitinated growth factor receptors, the beta-arrestins which bind GPCRs and ARH which binds LDL receptor family members. All of these function as cargo-selective clathrin adaptors, targeting surface receptors for internalization by clathrin-mediated endocytosis.

In beta-arrestin, the cargo motif is regulated by a remarkable structural rearrangement. The motif maintains the endocytosis-incompetent state by binding back on the folded core of the protein in a beta strand conformation. Apparently triggered via a beta-arrestin/GPCR interaction, the motif must be displaced and must undergo a strand to helix transition to enable the beta2 appendage binding step that drives GPCR-beta-arrestin complexes into clathrin coats.

The sandwich subdomain or side surface site binds an FxxxFxDF motif found in EPS15 and AP180 and may also accept other endocytosis proteins with variant motifs, as suggested by mutagenesis of the binding surface. The sandwich domain binders are accessory endocytic proteins (without a direct role in cargo binding) which help in CCV formation. Currently, there is no entry for the sandwich domain motif in ELM.

Figure 3. ELM detail page showing information about the ELM class TRG_AP2beta_CARGO_1.

Table 3. Main cellular compartments used in ELM annotation

Count	GO Id	GO term
98	GO:0005829	Cytosol
69	GO:0005634	Nucleus
17	GO:0005576	Extracellular
12	GO:0005794	Golgi apparatus
10	GO:0005886	Plasma membrane
9	GO:0009898	Internal side of plasma membrane
9	GO:0005783	Endoplasmic reticulum
6	GO:0005739	Mitochondrion
5	GO:0005643	Nuclear pore
5	GO:0045334	Clathrin-coated endocytic vesicle

background coloring to highlight structured regions within the protein, which allows inspection of SLIMs that reside at domain boundaries and emphasizes motifs in disordered regions.

The conservation of linear motifs can help in assessing the functional relevance of putative instances, with functional instances showing higher overall sequence conservation than non-functional ones (21). Therefore, sequence conservation of the query protein is calculated using a tree-based conservation scoring method (22) and highlighted in the graphical output. Here, lighter shades of blue represent low conservation while dark blue shading corresponds to high-sequence conservation. The actual conservation score can be inspected by moving the mouse over the respective ELM instance (Figure 2).

The functionality of linear motifs can be modulated by modifications such as phosphorylation (23,24). To enable the user to investigate phosphorylation data in the context of putative linear motif instances, phosphorylation annotations from the Phospho.ELM resource (25) have been added to the graphical output (Figure 2, top row).

Search ELM Instances								
Search	ELMs	Instances	Candidates	Links	About	News	Help	Diseases
Full-Text Search (to show all instances, enter 'all' or '') <input type="text" value="ap2"/>								
Filter by instance Logic	true positive							
Filter by organism	Homo sapiens							
<input type="button" value="submit"/>	<input type="button" value="Reset"/>							
58 Instances for search term 'ap2': (click table headers for sorting)								
CLV								
LIG								
MOD								
TRG								
ELM identifier	Sequence	Start	End	Subsequence	Instance Logic	#Evidence	PDB	Organism
LIG_AP2alpha_1	AMPH_HUMAN	324	328	QENIISFEDNEVPEISVTT	true positive	1	1KY7	Homo sapiens (Human)
LIG_AP2alpha_1	AP180_HUMAN	440	444	VTAEVDLFGDAFAASPGEAP	true positive	0	---	Homo sapiens (Human)
LIG_AP2alpha_1	AP180_HUMAN	564	568	APPALDIFGDLFESTPEVAA	true positive	0	---	Homo sapiens (Human)
LIG_AP2alpha_1	AP180_HUMAN	642	646	SSGVIDLFGDAFGSSASEPQ	true positive	0	---	Homo sapiens (Human)
LIG_AP2alpha_1	BIN1_HUMAN	362	366	QEQLISLFEDIFVPEISVTT	true positive	0	---	Homo sapiens (Human)
LIG_AP2alpha_2	EP15R_HUMAN	599	601	RGSFGANDDPFKNKKALLFSN	true positive	1	---	Homo sapiens (Human)
LIG_AP2alpha_2	EP15R_HUMAN	618	620	NNTOELHPDPFOTEDPKSD	true positive	1	---	Homo sapiens (Human)
LIG_AP2alpha_2	EP15R_HUMAN	624	626	HPDPFOTEDPDKSDPFKGAD	true positive	1	---	Homo sapiens (Human)
LIG_AP2alpha_2	EP15R_HUMAN	629	631	QTEDPFKSKDPFKGADPFKGD	true positive	1	---	Homo sapiens (Human)
LIG_AP2alpha_2	EP15R_HUMAN	640	642	KGADPFKGDPEONDPAEQQ	true positive	1	---	Homo sapiens (Human)
LIG_AP2alpha_2	EP15R_HUMAN	645	647	FKGDPFQNDPFEAEQQTSTD	true positive	1	---	Homo sapiens (Human)
LIG_AP2alpha_2	EP15R_HUMAN	656	658	AEQQTSTDTPFGDPFKESD	true positive	1	---	Homo sapiens (Human)
LIG_AP2alpha_2	EP15R_HUMAN	661	663	TSTDPFGGDPFKESDPFRGS	true positive	1	---	Homo sapiens (Human)
LIG_AP2alpha_2	EP15R_HUMAN	667	669	GGDPFKESDPFRGSATDDFF	true positive	1	---	Homo sapiens (Human)
LIG_AP2alpha_2	EP15R_HUMAN	685	687	FFKKDTKIDDPFTSDPFTKIP	true positive	1	---	Homo sapiens (Human)
LIG_AP2alpha_2	EP15R_HUMAN	690	692	TKNDPFTSDPFTKNPSPLSK	true positive	1	---	Homo sapiens (Human)
LIG_AP2alpha_2	EP15R_HUMAN	703	705	NPSLPSKLDPFESSDPFSSS	true positive	1	---	Homo sapiens (Human)
LIG_AP2alpha_2	EP15R_HUMAN	709	711	KLDPEFSSDPFESSSVSSKG	true positive	1	---	Homo sapiens (Human)
LIG_AP2alpha_2	EP15R_HUMAN	722	724	SSVSSKGSDFPFGTLDPPFGSG	true positive	1	---	Homo sapiens (Human)
LIG_AP2alpha_2	EP15R_HUMAN	728	730	GSDPFGTLDPFGSGFSINSAE	true positive	1	---	Homo sapiens (Human)
LIG_AP2alpha_2	EP15R_HUMAN	768	770	LGGAGFSDDPFKSKQDTPAL	true positive	1	---	Homo sapiens (Human)
LIG_AP2alpha_2	EP15R_HUMAN	812	814	SADFPEAPDPFOLGADSGD	true positive	1	---	Homo sapiens (Human)
LIG_AP2alpha_2	EP15R_HUMAN	823	825	QPLGADSDPFEOSKKGFDGP	true positive	1	---	Homo sapiens (Human)
LIG_AP2alpha_2	EP15R_HUMAN	833	835	F0SKKGFDGPFGSKDQFPVPS	true positive	1	---	Homo sapiens (Human)
LIG_AP2alpha_2	EP15R_HUMAN	839	841	FGDPFGSKDPEVPSSAKPS	true positive	1	---	Homo sapiens (Human)
LIG_AP2alpha_2	EPS15_HUMAN	599	601	NNRHSKEDDPFIVDSSSLTG	true positive	1	---	Homo sapiens (Human)
LIG_AP2alpha_2	EPS15_HUMAN	606	608	---	---	1	---	Homo sapiens (Human)

Figure 4. ELM instances browse page. A full-text search (here, search term used was ‘AP2’, filtering for ‘true positive’ instances in taxon ‘*Homo sapiens*’, yielding 58 instances) assists in finding annotated instances. A search can be restricted to a particular taxonomy or instance logic (top) or ELM class type (buttons on the left). The list can also be exported to TSV or FASTA format for further processing.

The phosphorylated residues are highlighted in different colors (serine: green, threonine: blue, tyrosine: red); each phosphorylation site is linked to a page showing detailed information about the respective modification site from the manually curated data set of the Phospho.ELM resource.

VIRAL INSTANCES

The importance of the short linear motifs in virus–host interactions makes the ELM resource an important tool for the viral research community. For example, Cruz *et al.* (26) analyzed a protein phosphatase 1 (PP1) docking motif in ‘protein 7’ of transmissible gastroenteritis virus using the ELM class LIG_PP1. This conserved sequence motif mediates binding to the PP1 catalytic subunit, a key

regulator of the cellular antiviral defense mechanisms, and is also found in other viral proteomes, suggesting that it might be a recurring strategy to counteract the hosts' defense against RNA viruses by dephosphorylating eukaryotic translation initiation factor 2 α and ultimately ribonuclease L.

To reflect our increasing awareness of viral motifs (8), special focus has been attributed to the annotation of viral instances in the ELM database: in the latest release, more than 200 novel ELM instances found in 84 different viral taxons have been added. The notion of viruses abusing existing SLiMs in their hosts is demonstrated by viral instances being annotated alongside instances in their hosts' proteins. For example, the ELM class LIG_PDZ_Class_1 contains 12 instances in human proteins but has recently been expanded with 5 instances from 5 different human pathogenic virus proteins.

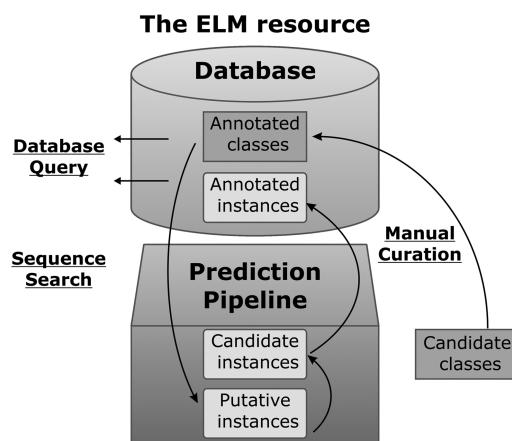


Figure 5. Schema of the ELM resource and data life cycle. Annotated ELM classes, and instances thereof, can be searched by database query. Via sequence search by the motif detection pipeline, annotated ELM classes yield putative instances in query sequences. By adding experimental evidence and references, these putative instances become candidate instances for annotation, and, with further curation, ultimately become fully annotated instances.

LINEAR MOTIFS AND DISEASES

The importance of SLiMs is further corroborated by the occurrence of pathologies that are caused by mutations that either mutate existing linear motifs or create novel linear motifs (of undesired function) (27). Examples include ‘Usher’s syndrome’ (28), ‘Liddle’s Syndrome’ (6) or ‘Golabi-Ito-Hall Syndrome’ (29). The developmental disorder ‘Noonan Syndrome’ can be caused by mutations in Raf-1 that abrogate the interaction with 14-3-3 proteins mediated by corresponding SLiMs and thereby deregulate the Raf-1 kinase activity (30) (the Raf-1 protein sequence features two LIG_14-3-3_1 binding sites that are annotated at 256–261 and 618–623 in the ELM resource). A related disease, ‘Noonan-like Syndrome’, is caused by an S to G mutation at position 2 of the SHOC2 protein, creating a novel myristylation site (annotated as ELM class MOD_NMyristoyl). This irreversible modification results in aberrant targeting of SHOC2 to the plasma membrane and impaired translocation to the nucleus upon growth factor stimulation (5). More information about the implication of short linear motifs on diseases is collected at <http://elm.eu.org/infos/diseases.html>.

APPLICATION OF THE ELM RESOURCE

By providing a high-quality, manually curated data set of linear motif classes with experimentally validated SLiM instances, the ELM database has proven to be invaluable to the community: small-scale (single protein) analyses benefit from the detailed annotation of each ELM class in attributing novel features to proteins of interest. By using *in vitro* and *in vivo* studies, von Nanelstadh *et al.* (31) could validate a PDZ class III motif, detected by ELM at the carboxy terminus of myotilin and the FATZ (calsarcin/myozenin) families. This evolutionarily conserved carboxy-terminal motif mediates binding to PDZ domains of ZASP/Cypher and other Enigma family members (ALP, CLP-36 and RIL) and disruption

of these interactions results in myofibrillar myopathies (32). Additionally, ELM annotations can contribute to high-throughput screenings (33) as well as development of novel algorithms (34–36), methods (37) and databases (38). Furthermore, the highly curated data of the ELM resource are used as a benchmarking data set to evaluate the accuracy of prediction algorithms (21,39,40).

For any such analysis, the user should be aware that many matches to ELM regular expressions are false positives. Before conducting experiments based on ELM results, it is strongly advisable to check if a motif match is conserved, exposed in a cell compartment in which the motif is known to be functional. The ELM resource applies several filters to provide the user with such information that should ideally also be supported by the experimental evidence.

SUMMARY

The importance of SLiMs is highlighted by the growing number of instances with relevance to diseases or viruses. Yet, despite their importance and abundance, our understanding of linear motifs is still limited. This is mainly owing to the fact that they are still quite difficult to predict computationally and to investigate experimentally (3,41,42). By better understanding the biology of linear motifs, we hope to increase our insight into diseases and viruses (and vice versa). The ELM resource tries to aid the researcher in the search for putative SLiM instances by providing a feature-rich toolset for sequence analysis. Consequently, with the aforementioned additions and changes, we hope that the ELM resource continues to be a valuable asset to the community.

ACKNOWLEDGEMENTS

The authors would like to thank the users of the ELM resource as well as all colleagues, contributors and annotators of the ELM resource.

FUNDING

EMBL international PhD program (to R.J.W.); EMBL Interdisciplinary PostDoc fellowship (EIPOD to N.E.D.); NGFN framework by the Federal Government Department of Education and Science [FKZ01GS0862 (DiGToP) to M.S. and M.H.]; European Community’s Seventh Framework Programme FP7/2009 (SysCilia) (241955 to G.T.) and (SyBoSS) (242129 to K.V.R.); Polish Ministry of Science and Higher Education within Iuventus Plus project (IP2010-0483-70 to M.D.); Biotechnology and Biological Sciences Research Council (BB/F010486/1 to A.C.); Région Alsace and Collège Doctoral Européen (to K.L.); Science Foundation Ireland (08/IN.1/B1864 to G.G.); BBSRC New Investigator Award (BB/I006230/1 to R.J.E.); German Research Foundation (SFB796 Project A2 to H.M.); grants from the Swiss National Science Foundation (to M.O.S.). Funding for open access charge: EMBL.

Conflict of interest statement. None declared.

REFERENCES

- Davey,N.E., Van Roey,K., Weatheritt,R.J., Toedt,G., Uyar,B., Altenberg,B., Budd,A., Diella,F., Dinkel,H. and Gibson,T.J. (2011) Attributes of short linear motifs. *Mol Biosyst.*, September 12 (doi:10.1039/c1mb05231d; epub ahead of print).
- Fuxreiter,M., Tompa,P. and Simon,I. (2007) Local structural disorder imparts plasticity on linear motifs. *Bioinformatics*, **8**, 950–956.
- Diella,F., Haslam,N., Chica,C., Budd,A., Michael,S., Brown,N.P., Trave,G. and Gibson,T.J. (2008) Understanding eukaryotic linear motifs and their role in cell signaling and regulation. *Front. Biosci.*, 6580–6603.
- Gibson,T.J. (2009) Cell regulation: determined to signal discrete cooperation. *Trends Biochem. Sci.*, **34**, 471–482.
- Cordeddu,V., Di Schiavi,E., Pennacchio,L.A., Ma'ayan,A., Sarkozy,A., Fodale,V., Cecchetti,S., Cardinale,A., Martin,J., Schackwitz,W. et al. (2009) Mutation of SHOC2 promotes aberrant protein N-myristylation and causes Noonan-like syndrome with loose anagen hair. *Nat. Genet.*, **9**, 1022–1026.
- Furuhashi,M., Kitamura,K., Adachi,M., Miyoshi,T., Wakida,N., Ura,N., Shikano,Y., Shinshi,Y., Sakamoto,K., Hayashi,M. et al. (2005) Liddle's syndrome caused by a novel mutation in the proline-rich PY motif of the epithelial sodium channel beta-subunit. *J. Clin. Endocrinol. Metab.*, **1**, 340–344.
- Deretic,D., Schmerl,S., Hargrave,P.A., Arendt,A. and McDowell,J.H. (1998) Regulation of sorting and post-Golgi trafficking of rhodopsin by its C-terminal sequence QVS(A)PA. *Proc. Natl Acad. Sci. USA*, **18**, 10620–10625.
- Davey,N.E., Trave,G. and Gibson,T.J. (2011) How viruses hijack cell regulation. *Trends Biochem. Sci.*, **3**, 159–169.
- Neduva,V., Linding,R., Su-Angrand,I., Stark,A., de Masi,F., Gibson,T.J., Lewis,J., Serrano,L. and Russell,R.B. (2005) Systematic discovery of new recognition peptides mediating protein interaction networks. *PLoS Biol.*, **12**, e405.
- Hulo,N., Bairoch,A., Bulliard,V., Cerutti,L., Cuche,B.A., de Castro,E., Lachaize,C., Langendijk-Genevaux,P.S. and Sigrist,C.J. (2008) The 20 years of PROSITE. *Nucleic Acids Res.*, **36**, D245–D249.
- Rajasekaran,S., Balla,S., Gradie,P., Gryk,M.R., Kadaveru,K., Kundeti,V., Maciejewski,M.W., Mi,T., Rubino,N., Vyas,J. et al. (2009) Minimotif miner 2nd release: a database and web system for motif search. *Nucleic Acids Res.*, **37**, D185–D190.
- Obenauer,J.C., Cantley,L.C. and Yaffe,M.B. (2003) Scansite 2.0: proteome-wide prediction of cell signaling interactions using short sequence motifs. *Nucleic Acids Res.*, **13**, 3635–3641.
- Puntervoll,P., Linding,R., Gemund,C., Chabanis-Davidson,S., Mattingdal,M., Cameron,S., Martin,D.M., Ausiello,G., Brannetti,B., Costantini,A. et al. (2003) ELM server: a new resource for investigating short functional sites in modular eukaryotic proteins. *Nucleic Acids Res.*, **13**, 3625–3630.
- Gould,C.M., Diella,F., Via,A., Puntervoll,P., Gemund,C., Chabanis-Davidson,S., Michael,S., Sayadi,A., Bryne,J.C. et al. (2010) ELM: the status of the 2010 eukaryotic linear motif resource. *Nucleic Acids Res.*, **38**, D167–D180.
- Davey,N.E., Edwards,R.J. and Shields,D.C. (2010) Computational identification and analysis of protein short linear motifs. *Front. Biosci.*, **15**, 801–825.
- Hermjakob,H., Montecchi-Palazzi,L., Bader,G., Wojcik,J., Salwinski,L., Ceol,A., Moore,S., Orchard,S., Sarkans,U., von Mering,C. et al. (2004) The HUPO PSI's molecular interaction format—a community standard for the representation of protein interaction data. *Nat. Biotechnol.*, **2**, 177–183.
- Cote,R.G., Jones,P., Apweiler,R. and Hermjakob,H. (2006) The Ontology Lookup Service, a lightweight cross-platform tool for controlled vocabulary queries. *BMC Bioinformatics*, **7**.
- Via,A., Gould,C.M., Gemund,C., Gibson,T.J. and Helmer-Citterich,M. (2009) A structure filter for the Eukaryotic Linear Motif Resource. *BMC Bioinformatics*, **10**.
- Letunic,I., Doerks,T. and Bork,P. (2009) SMART 6: recent updates and new developments. *Nucleic Acids Res.*, **37**, D229–D232.
- Dosztanyi,Z., Csizmok,V., Tompa,P. and Simon,I. (2005) IUPred: web server for the prediction of intrinsically unstructured regions of proteins based on estimated energy content. *Bioinformatics*, **16**, 3433–3434.
- Dinkel,H. and Sticht,H. (2007) A computational strategy for the prediction of functional linear peptide motifs in proteins. *Bioinformatics*, **24**, 3297–3303.
- Chica,C., Labarga,A., Gould,C.M., Lopez,R. and Gibson,T.J. (2008) A tree-based conservation scoring method for short linear motifs in multiple alignments of protein sequences. *BMC Bioinformatics*, **9**.
- Balagopalan,L., Coussens,N.P., Sherman,E., Samelson,L.E. and Sommers,C.L. (2010) The LAT story: a tale of cooperativity, coordination, and choreography. *Cold Spring Harb. Perspect. Biol.*, **8**, a005512.
- Pawson,T. and Scott,J.D. (2005) Protein phosphorylation in signaling—50 years and counting. *Trends Biochem. Sci.*, **6**, 286–290.
- Dinkel,H., Chica,C., Via,A., Gould,C.M., Jensen,L.J., Gibson,T.J. and Diella,F. (2011) Phospho.ELM: a database of phosphorylation sites—update 2011. *Nucleic Acids Res.*, **39**, D261–D267.
- Cruz,J.L., Sola,I., Becares,M., Alberca,B., Plana,J., Enjuanes,L. and Zuniga,S. (2011) Coronavirus gene 7 counteracts host defenses and modulates virus virulence. *PLoS Pathog.*, **6**, e1002090.
- Kadaveru,K., Vyas,J. and Schiller,M.R. (2008) Viral infection and human disease—insights from minimotifs. *Front. Biosci.*, **13**, 6455–6471.
- Weil,D., El-Amraoui,A., Masmoudi,S., Mustapha,M., Kikkawa,Y., Laine,S., Delmaghani,S., Adato,A., Nadifi,S., Zina,Z.B. et al. (2003) Usher syndrome type I G (USH1G) is caused by mutations in the gene encoding SANS, a protein that associates with the USH1C protein, harmonin. *Hum. Mol. Genet.*, **5**, 463–471.
- Tapia,V.E., Nicolaescu,E., McDonald,C.B., Musi,V., Oka,T., Inayoshi,Y., Satteson,A.C., Mazack,V., Humbert,J., Gaffney,C.J. et al. (2010) Y65C missense mutation in the WW domain of the Golabi-Ito-Hall syndrome protein PQBP1 affects its binding activity and deregulates pre-mRNA splicing. *J. Biol. Chem.*, **25**, 19391–19401.
- Pandit,B., Sarkozy,A., Pennacchio,L.A., Carta,C., Oishi,K., Martinelli,S., Pogna,E.A., Schackwitz,W., Ustaszewska,A., Landstrom,A. et al. (2007) Gain-of-function RAF1 mutations cause Noonan and LEOPARD syndromes with hypertrophic cardiomyopathy. *Nat. Genet.*, **8**, 1007–1012.
- von Nandelstadh,P., Ismail,M., Gardin,C., Suila,H., Zara,I., Belgrano,A., Valle,G., Carpen,O. and Faulkner,G. (2009) A class III PDZ binding motif in the myotilin and FATZ families binds enigma family proteins: a common link for Z-disc myopathies. *Mol. Cell Biol.*, **3**, 822–834.
- Selcen,D. and Engel,A.G. (2004) Mutations in myotilin cause myofibrillar myopathy. *Neurology*, **8**, 1363–1371.
- Gfeller,D., Butty,F., Wierzbicka,M., Verschueren,E., Vanhee,P., Huang,H., Ernst,A., Dar,N., Stagljar,I., Serrano,L. et al. (2011) The multiple-specificity landscape of modular peptide recognition domains. *Mol. Syst. Biol.*, **7**.
- Bauer,D.C., Willadsen,K., Buske,F.A., Le Cao,K.A., Bailey,T.L., Dellaire,G. and Boden,M. (2011) Sorting the nuclear proteome. *Bioinformatics*, **13**, i7–i14.
- Walsh,I., Martin,A.J., Di Domenico,T., Vullo,A., Pollastri,G. and Tosatto,S.C. (2011) CSpritz: accurate prediction of protein disorder segments with annotation for homology, secondary structure and linear motifs. *Nucleic Acids Res.*, **39**, W190–W196.
- Lieber,D.S., Elemento,O. and Tavazoie,S. (2010) Large-scale discovery and characterization of protein regulatory motifs in eukaryotes. *PLoS One*, **12**, e14444.
- Pless,O., Kowen-Leutz,E., Dittmar,G. and Leutz,A. (2011) A differential proteome screening system for post-translational modification-dependent transcription factor interactions. *Nat. Protoc.*, **3**, 359–364.

38. Goel,R., Muthusamy,B., Pandey,A. and Prasad,T.S. (2011) Human protein reference database and human proteinpedia as discovery resources for molecular biotechnology. *Mol. Biotechnol.*, **1**, 87–95.
39. Edwards,R.J., Davey,N.E. and Shields,D.C. (2007) SLiMFinder: a probabilistic method for identifying over-represented, convergently evolved, short linear motifs in proteins. *PLoS One*, **10**, e967.
40. Edwards,R.J., Davey,N.E. and Shields,D.C. (2008) CompairMotif: quick and easy comparisons of sequence motifs. *Bioinformatics*, **10**, 1307–1309.
41. Perkins,J.R., Diboun,I., Dessailly,B.H., Lees,J.G. and Orengo,C. (2010) Transient protein-protein interactions: structural, functional, and network properties. *Structure*, **10**, 1233–1243.
42. Edwards,R.J., Davey,N.E., Brien,K.O. and Shields,D.C. (2011) Interactome-wide prediction of short, disordered protein interaction motifs in humans. *Mol. Biosyst.*, August 30 (doi:10.1039/c1mb05212h; epub ahead of print).